



Espacenet

Bibliographic data: JP6204952 (A) — 1994-07-22

TRAINING OF SPEECH RECOGNITION SYSTEM UTILIZING TELEPHONE LINE

Inventor(s): BINSU EMU SUTANFUOODO; NOOMAN EFU
BURITSUKUMAN ±

Applicant(s): IBM ±

Classification: - **international:** G06F3/16; G10L15/00; G10L15/06;
H04B14/04; G10L15/02; G10L15/14;
G10L15/20; G10L19/00; (IPC1-
7): G06F3/16; G10L5/06; H04B14/04
- **european:** G10L15/06T

Application number: JP19930219208 19930812

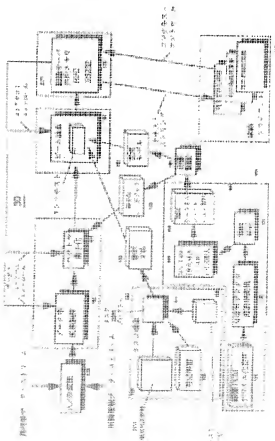
Priority number(s): US19920948031 19920921

Also published as: JP2524472 (B2) US5475792 (A)

Abstract of JP6204952 (A)

PURPOSE: To provide a system with which continuous voices inputted from a telephone line can be recognized by inputting a correction voice data set to a voice recognizing processor and training a statistical pattern matching unit.

CONSTITUTION: A digital voice input from an A/D 100 is sampled by a data rate converting device 102, a voice data stream is divided into fine frames by a vector-quantizing block 104, and characteristics are extracted from the respective frames and summarized. Vector quantization is similarly used in a training process 190, as well. A quantizing code book 105 preserves a code book for telephone voice generated by the function 190. On the other hand, input voice data 180 are resampled in a block 182, supplied to a code/decode band filter and set to a prescribed band width. The correction voice data set is inputted to the voice-recognizing processor, and the statistical pattern matching unit is trained. The voice recognition is executed, while using a voice recognition processor from a telephone system.



Last updated
5/12/2011 Worldwide Database 5/7/31
93p

1

【特許請求の範囲】

【請求項1】 電話帯域幅より高い帯域幅の音声認識訓練プロセッサへの音声データ・セットの入力ステップと、

上記音声データ・セットを間引き、上記電話帯域幅を有する間引かれた音声データ・セットを入手するステップと、

帯域通過デジタル濾波器を上記間引かれた音声データ・セットに適用し、電話機種の電送特性に特化した、濾波された音声データ・セットを入手するステップと、

上記濾波された音声データ・セットの振幅を、その最大ダイナミック・レンジが非圧伸電話音声の最大ダイナミック・レンジと一致するように補正し、振幅補正音声データ・セットを入手するステップと、

上記振幅補正音声データ・セットを、電話システムの圧伸・非圧伸音声信号シークエンスを表す量子化ノイズを用いて修正し、修正音声データ・セットを入手するステップと、

上記修正音声データ・セットを音声認識プロセッサに入力し、統計的パターン・マッチング・ユニットを訓練するステップと、

から構成される、電話システムから得られる音声に 대응する音声認識プロセッサを訓練する方法。

【請求項2】 上記電話帯域幅が上記音声幅の高位帯域より低い帯域である上記請求項1記載の方法。

【請求項3】 上記帯域通過デジタル濾波器が最大平坦設計アルゴリズムを備え持つ上記請求項1記載の方法。

【請求項4】 上記音声データ・セット振幅補正の結果、最大ダイナミック・レンジが非圧伸 μ -law 電話音声の最大ダイナミック・レンジに一致する上記請求項1記載の方法。

【請求項5】 上記音声データ・セット振幅補正の結果、最大ダイナミック・レンジが非圧伸 μ -law 電話音声の最大ダイナミック・レンジに一致する上記請求項1記載の方法。

【請求項6】 上記音声データ・セット修正ステップが μ -law ノイズとしての量子化ノイズを用いる上記請求項1記載の方法。

【請求項7】 上記音声データ・セット修正ステップが μ -law ノイズとしての量子化ノイズを用いる上記請求項1記載の方法。

【発明の詳細な説明】

【0001】

【産業上の利用分野】本発明は、公衆電話交換回線を利用する音声認識システムに関するものである。

【0002】

【従来の技術】音声認識システムは、よく知られている技術である。IBM タングラ (Tangora) [13] (本願書文末記載の参照文献の番号で、以下同様に表記する) およびドラゴン・システム・ドラゴン 30 k 口述システムは

2

その例である。それらは、典型的な単一ユーザおよび話し手依存型システムである。これは、プロセスが「登録」と呼ばれるプロセスの間に、話し手の音声パターンで音声認識装置を訓練することを各話し手に要求する。

将来の認識セッションの中で話し手自身をシステムが識別しなければならないのでシステムは話し手のプロフィールを維持する。典型的には、話し手は低レベル雑音システム環境の中でローカル・マイクを通して認識を行う。将来の認識セッションの中で話し手自身をシステムが識別しなければならないのでシステムは話し手のプロフィールを維持する。典型的には、話し手は低レベル雑音システム環境の中でローカル・マイクを通して認識を行う。

登録作業の間、その話し手は、長つたらしい原稿を読むことを要求されるが、それ故に、そのシステムは各話し手の特色に順応することができることとなる。独立した口述システム (たとえば、上記の2つのシステム) は、話し手にたどらざる、不自然な形で、すなわち、語と語の間にポーズをいれながら、各語を形づくることを要求する。これにより、音声認識システムは、語の境界となる、先行および後続の無音を利用し、各個人の語に連想される音声パターンを識別することが可能となる。

典型的音声認識システムは、(たとえば、IBM タングラ システムの Office Correspondence の場合のように) 単一の機械上で動作し、訓練された単一の適用業務を持つ。

【0003】話し手依存型音声認識装置をもつマルチ・ユーザ・システム環境は、各話し手にその音声パターンをシステムに理解させるための退屈な訓練に従事することを要求する。話し手の電話番号によってシステムがどの音声テンプレートを使用すべきかを知り得る共通データ・ベースに音声テンプレートが格納されているかもしれないが、それでもなお各話し手は使用の前にそのシステムを訓練しなければならない。外部の電話線から接続してくる新しいシステム利用者は、このプロセスが容認できるものでないことを認識する。また、成功した電話の音声認識システムというものは、様々な分野に関係する音声を正確に認識するために迅速な文脈切り替えができなければならない。たとえば、一般のオフィス通信のために訓練されたシステムは、数字列の提示の場合、うまく働かない。

【0004】Kai-Fu Lee の博士号論文 [1] の中で最初に記述された スフィンクス (Sphinx) システムは、以前の話し手依存型認識システムに大きな進歩をもたらした。それは話し手独立型であり、会話音声の連続ストリームから単語を認識することができた。このシステムは、実際の使用に先立って行われる話し手個々の登録を必要としなかった。話し手依存型システムの中には、話し手に4〜6週毎に再登録することを要求したり、利用者にそのシステムが理解するための個人用アラグレン・カードリッジを持ち運ぶことを要求する。連続の音声認識を行うスフィンクス・システムは、語と語の間の休止を必要とせず、音声認識システムの一時的ユーザに非常に多くの微細なアプローチを提供する。認識システムの利便的な

めにどのように音声を調節するかをユーザが訓練しなくてよいので、この点は、電話の音声認識システムの本質的な特長である。

【0005】音声認識システムは、また、与えられたさまざまな語彙を使って、実時間処理を提供しなければならない。しかし、スフィンクス・システムは、まだ以前の話し手に依存する認識システムの不利な点をいくつか持っていた。マイクロホンおよび比較的制約された語彙を使用しながら低レベル雑音システム環境の中で単一機械上で操作するようプログラムされていた。スフィンクス・システムは、複数ユーザのサポート、少くとも、異なるロケーションおよび複数の語彙認識に関するサポートを行うようには設計されなかった。

【0006】

【発明が解決しようとしている課題】本発明は、上記の従前技術の不利益の多くを克服することを目的とする。したがって、本発明は、ローカルおよび遠隔地及方の話し手からの入力を持つ電話機器使用に適した連続音声話し手独立型音声認識システムを提供することを目指す。

【0007】低レベル雑音条件の中で集められた語彙を基に電話システム環境のように高レベル雑音の中での音声パターンを認識できるようにシステムを訓練することは、本発明のもうひとつの目的である。

【0008】複数の音声適用業務が、コンピュータ・ネットワーク上または電話線上で同時に音声認識システムによって音声認識できるようにすることは、本発明のもうひとつの目的である。

【0009】

【課題を解決するための手段】本発明の上記目的は、ローカル・エリア・ネットワークまたは広域ネットワークの上のクライアント・サーバを基に構築される音声認識システムによって達成される。この音声認識システムは、アナログまたはデジタル音声データを音声を表わす一組のケプストラム係数およびベクトル量子化値に変換するフロントエンドを含む多くのモジュールに分けられる。バックエンドは、ベクトル量子化値を使用して、その音声の作る文脈と言葉モデル(Phoneme Models)と語対文法(Word Pair Grammars)に従ってその語を認識する。語彙を一連の文脈(すなわち、ある特定の語がそのシステムによって予期される状況)に分割することによって、一層大規模な語彙を、最小限のメモリに収納することができる。ユーザが音声認識作業を進めるにつれて、文脈は共通のデータベースから迅速に切り換えられる(下記引用Brickmanその他による特許出願参照)。システムは、また、コンピュータ・ネットワーク間および複数のユーザ適用業務間のインターフェースを備える。

【0010】このシステムは、文脈のための語対文法を構築しシステムを訓練する訓練およびタスク構築モジュールを備え持つ。

【0011】本発明は、電話から得られる音声に 대응できるように音声認識システムを訓練するための電話チャネル・シミュレーション・プロセスを含む。その方法は、音声データ・セットを、電話帯域幅より帯域幅が高い音声認識訓練・プロセッサに入力することから始める。入力音声データ・セットは、電話帯域幅を持つ間引かれた(decimated)音声データを得るために間引かれる。その後、帯域通過デジタル濾波器(Bandpass Digital Filter)を間引かれた音声データ・セットに適用し、電話機器の伝送特性に特化させる。これにより、濾波された音声データ・セットが得られる。次に、その濾波された音声データ・セットの最大ダイナミック・レンジが非圧伸(uncompanded)電話音声の最大レンジと一致するように、振幅補正(スケージング)を行う。それから、振幅補正された音声データ・セットは、電話システムの音声信号を圧伸、非圧伸するシーケンスを表わす量子化ノイズによって修正される。そして修正された音声データは、統計的パターン・マッチング・データ・ユニットを訓練するために音声認識プロセッサに入力される。上記方法により、音声認識プロセッサは電話システムからの音声信号に対して音声認識を実行することができることとなる。

【0012】

【実施例】電話ラインによってもたらされる帯域幅減衰および雑音は、すべての音声認識システムの正確度を減ずる。この影響は、瞬間的に認識されなければならない語彙の大きさに応じて増加する。迅速に切り替え可能な音声認識文脈の使用はこの発明にとって有用であるが、そのために、個々の文脈のサイズが制限されなければならない。文脈切り替えは、N.F. Brickmanその他の発明になるIBM出願のアメリカ合衆国特許出願番号947,634 "Instantaneous Context Switching For Speech Recognition Systems"で記述されており、本明細書においても参照される。図1は、ハードウェア機械構成から独立したIBM連続音声認識システム(IBM Continuous Speech Recognition System)のことで、以下ICSRSSと呼ぶ)の論理的構造を図示する。ICSRSSは、幅広いレベルで、以下の分野をカバーするコンポーネントから構成される。

【0013】データ収集：データは、アナログからデジタル形式にブロック100で変換されるか、あるいは電話のデータの場合他のチャネルから潜在的にデマルチプレックス(demultiplexed)される。

【0014】データ圧縮：ICSRSSフロントエンドブロック102および104は、ベクトル量子化ステップの間に300バイト/秒に音声データストリームを調整し、再標準化し、圧縮する。

【0015】音声認識：バックエンド106は、文法ガイド型ビーム・サーチ・アルゴリズムを使用しているパターンマッチング音楽モデル192によって実際の音声

5

認識を実行する。音楽モデル192および語対文法135は共に認識文脈を構成する。バックエンド認識装置のひとつまたは複数の事例が、遠隔地であろうがローカルであろうが音声データを捕捉して圧縮するフロントエンド事例に配備されることである。

【0016】タスク構築：タスク構築コンポーネント130は、認識文脈のオフラインでの構築を可能にし、実行時で使用のために語対文法をコンパイルし、適切な音楽モデルをそのタスク（文脈）に連結させる。

【0017】適用業務プログラム・インタフェース（API）：API108は、データストリーム・コントロール、文脈ローディングおよび起動を可能にするRPC（Remote Procedure Call）に基づく認識サービスを提供する。

【0018】電話チャネル・シミュレータ：シミュレータ185は、高帯域、高解像度音声データ・セットを、音楽モデル192および電話音声に連結し、減少された標本抽出率、圧縮された帯域幅および圧縮されたダイナミック・レンジの電話音声を作り出す。

【0019】音声認識の間に、ローカル・マイクからの高帯域音声データストリームも電話に関連しているような低帯域音声データストリームも、アナログデジタル変換ブロック100によって受け取られる。アナログデジタル変換100は、ボイス・ワークステーション上のIBM M-Audio Capture/Playback Cardカード（M-ACPA）のようなハードウェア・カードによって実行されることができ、M-ACPAは、高帯域または電話帯域幅信号を処理するデジタル信号処理機構を持ち、デジタルに標本化された一連のデータ・ポイントにそれらを変換する。この変換は、また、デジタルPBXや8KHz、8ビットのMu-Law/A-Law形式で与えられる電話データストリームによって実行されることもできる。

【0020】本発明では、高帯域を、サンプリング率16キロヘルツ以上と定義する。低レベル帯域幅を、アメリカ合衆国で一般の電話がデジタル音声に使う8キロヘルツ以下と定義する。電話システムの中でデジタル情報が個人の電話交換（PBX）から入る可能性があるため、A/D変換ブロック100は、オプションとして必要である。

【0021】音声認識に対する「フロントエンド」の中の最初の重要なブロックは、データ条件付け・速度変換ブロック102（Data Conditioning and Rate Conversionのこと、以下DCRCと呼ぶ）である。A/D変換100からのデジタル化された入力は、44または8KHzである。本発明で間引き（DECIMATION）と呼び使用する再標本化テクニックは、IEEEの文献(2)によって提供されている。DCRC102は、デジタル化された信号に対しアンチエイリアシング（Anti-aliasing）・フィルターを使用し標本化を行い、次のステップでの使

6

用のために、16KHzまたは8KHzデータストリームを作る。DCRCおよびベクトル量子化プロセスは、以下に詳細に記述される。

【0022】音声認識の中でデータ条件付け・速度変換の後、音声データは、ベクトル量子化ブロック104に渡される。ベクトル量子化の中でデジタル・データ・ストリームは、1秒間の1/50のフレームに細分化され、16KHz、11KHzおよび8KHzそれぞれ標本化率率に対し各々320個、220個および160個の標本となる。本発明の好ましい実施例のひとつでは、いかなる帯域幅音声信号からも計算される1秒につき100フレームがあり、それらは50パーセント上重ねられ、ハミング・ウィンドウ（Hamming Window）が適用される。ハミング・ウィンドウは、文献(3)で定義されている。

【0023】音声データストリームがフレームに細分化されたあと、ベクトル量子化ステップは、各フレームから特性を抽出する。ベクトル量子化ステップの抽出部分で、LPCケプストラム係数と呼ばれる一連のパラメータが、計算される。ケプストラム係数は、パターン認識のために音声の重要な特性のいくつかを抜き出し、要約する。データの各フレームの中で、音声の1秒の50分の1が、ケプセルに入れられる。1秒につき50のフレームと想定するであろうが、50パーセントの上重ねがあるので、1秒につき100フレームが生成される。ケプストラム係数を計算するために、まず（ \cos ine bellである）ハミング・ウィンドウが、音声データに適用される。抽出されたデータが、無限時間連続フーリエ変換にあるようにするために、ハミング・ウィンドウは、音声データの各フレームのエッジを次第に減少させる。

【0024】ハミング・ウィンドウ化されたフレームは音声スペクトルを平坦にするために、そのZ変換が $1.0 - 0.97e^{-j\omega}$ （[149ページ参照]）であるところの濾波器を使用して事前に濾波される。それから、14個の自己相関係数が計算される。自己相関係数が、文献(4)の記述でよく知られている方法でケプストラム係数を計算するために使われる。13個のケプストラム係数は、14個の自己相関係数から引き出される。自己相関係数の数やケプストラム係数の次元数を変えることは可能である。これらの係数の統計的特性は、最終的なベクトル量子化ステップをガイドするために使われる。

【0025】ベクトル量子化は、訓練プロセス190の中でも同様に使われる。下記の訓練データの調整は、基本スフィックス認識エンジンで電話機器上で動作可能とされる点で、本発明にとって重要である。訓練プロセス190において、10,000から15,000の間のセンテンスがとられて、フレームに細分化され、そこから自己相関およびケプストラム係数が計算される。参照文献(5)に記述されるK-手法タイプのクラスタリング

7

・プロシーダを使用して、256個のクラスにケプストラム・フレーム特性を区分する。これらのケプストラム・クラスターの中央値、およびそのクラス・ラベルが共に取り出され、これ以後「コード・ブック」と呼ばれる。量子化コード・ブック105は、音響訓練機能190によって生成される電話音声用コード・ブックを保存し、また、第2の高帯域音声用コード・ブックをも保存する。

【0026】ベクトル量子化の最終的なステップのために、どのクラスター中央値がフレーム・ケプストラム係数に最も近いかを決定するために、ブロック104は、上記のように訓練プロシーダで引き出される量子化コード・ブック105のコード・ブックを参照する。現在のフレームが、コード・ブック値によって表わされたクラスに割り当てられる。256個のクラスがあるので、VQ (Vector Quantization) 値は、1バイトで表わされる。微分ケプストラムおよびフレームのそのべき乗から引き出される別の2個の1バイトVQ値がある。1秒に100回引き出される3個の1バイトVQ値があり、その結果、音声データストリームは2、400 ビット/秒に圧縮される。

【0027】音声認識装置のためにその音声の特徴づける点の完全な別個のコード・ブックが、電話データから引き出され、図1の量子化コード・ブック105で保存されなければならないということは、電話音声認識に関する本発明の一部である。また、対応する音楽モデルが電話データから引き出され、音楽モデル192で保存されなければならないということは、本発明のもう一つの部分である。標準率減少、帯域幅圧縮およびダイナミック・レンジ圧縮のために、電話音声信号はかなり変わる。しかし、多大な努力を要する、電話から収集する音声標本の使用を必要とせず、高帯域標本を、電話チャネル特性をシミュレートするように処理することができる。これにより、スフィックス・システムの初期化訓練で使われた、大規模で既に使用可能な音声データ・ファイルを活用して、電話音声認識を可能となる。電話チャネル・シミュレータは、本発明の対象である。

【0028】電話チャネル・シミュレーションは、下記の通り、3つの段階的プロセスで達成される。

1.) 電話帯域幅への変換

文献[14]から[19]で参照されるように、(たとえば44、100 Hz、あるいは16、000 Hzで集められた16ビット解像度データのような) 高帯域、高解像度音声データ・セットが図1のブロック180への入力となる。

【0029】入力音声データ・セット180は、最初に、図1のブロック182の中で[2]で記述の再標本化プログラムを使用して、8、000 Hzに再標本化される。このデータは、図1の機能ブロック182で、参考文献[8]で記述のMAXFLATルーチンの修正版を

8

使用して設計された符号化器帯域濾波器に供給される。この濾波器は、図2、3および4の中で図示される。この濾波器の通過帯域特性は、現代の米国における電話機の中で使われる符号化/復号化濾波器に近似するよう設計される。通過帯域、3dbポイントおよび移行(BANDWIDTH)帯域幅の設定は、本発明の有効性にとって重要である。ローカル電話回線上の音声に対する良好な認識を行う認識訓練のための符号化濾波器を設計するのは可能であるが、遠隔地の電話については難しい。そのような問題を避けるために、上記の特性は、たとえば、低位の3dbポイントに対しては300 Hz、上位の3dbポイントに対しては3、600 Hzに設定すべきである。移行帯域幅は、それぞれ、400 Hzおよび800 Hzでなければならない。通過帯域は500 Hzから3、200 Hzになる。実際の符号化器帯域濾波器の幅に近似するために、通過帯域リプルは、全通過帯域にわたり、1単位から0.1パーセント以上の偏差であってはならない。

【0030】スフィックス音声認識エンジンおよびタンブラーを始めとするその他の音声認識エンジンが線形濾波器によって提示されるスペクトルのひずみを感じ取る点に、注意することは重要である。スペクトルのひずみは、主要な音声認識特性(例えばケプストラム)が周波数スペクトルから引き出されるので、その通過帯域の中の平坦な周波数応答を持たない、複雑な認識作業については、いくぶん平坦な通過帯域応答からのマイナーな偏差が、本願発明者の研究室において観察され、結果として、絶対認識誤り率が数パーセント劣化した。したがって、最大平坦設計アルゴリズムは、必要である。「スペクトルの傾き」へのスフィックス音声認識エンジンの感度が、参考文献[9]の中で指摘された。したがって、MAXFLATまたは比較的低レベルの通過帯域のリップル設計は、必要とされる。

【0031】4、100 Hzから8、000 Hzへの再標本化率変換は、参考文献[8]の中で提供されたMAXFLATには過度な要求であり、それは、帯域通過特性が符号化器帯域濾波器に必要なとき、低通過帯域フィルターの設計のためにのみ役立てられる。このルーチンに対するデザイン特性は、0.5ヘマップするナイキスト周波数と1.0ヘマップする標本化周波数によって、正規化された周波数の3dbポイントおよび移行帯域幅を表わす2個のパラメータ、ベータおよびガンマによって与えられる。Kaiserの参考文献[8]によって、ガンマは「0.005より大きく小さい」値に制限されなければならないことが示唆されている。これより低い値では、使われる計算精度浮動小数点数を増やすためにルーチンの修正が要求であり、そのような濾波器の案件数は、およそガンマの2乗に反比例するので、フィルター係数バッファを200から4096に拡張する必要がある。このため、44、100 Hzから8、000 Hzへの交換に必要となる0.05の約10分の1または0.0

0.5のガンマ値をもつ低域フィルタと、2個の低域通過フィルタ設計、低域から高域通過帯域変換、および、低域と高域通過フィルタの組み合わせが、必要な帯域通過特性を実現するために要求された。

【0032】上記フィルタ設計の実現によって、4、4、100Hzデータは、参照文献(2)で記述される再標準アルゴリズムを使用して、図1の機能ブロック182の中で8、000Hzに変換され、米国長距離電話機器のための通過帯域に非常に近い符号化器通過帯域を提供する。このデータは、下記のステップ2および3に従って処理され、16ビットの、低雑音信号となる。

【0033】同様の通過帯域特性および速度低減削減は、この訓練テクニックの中で使われる16、000Hz音声サンプルのために必要であるが、例外は、移行バンド要求がそれほど要求していない点と低域加重が、要求された通過帯域平坦度特性を達成するにはそれほど必要とされない点とである。図2、3、4で、事前訓練再標準化操作の訓練に実行されたのと同様に、符号化器通過帯域のインパルス(Impulse)、マグニチュード(Magnitude)およびログ・マグニチュード(Log Magnitude)応答を再びを示す。

2) ダイナミック・レンジを正規化するための振幅補正(スケール)

音声標準は、個別に読まれて、図1のブロック184で、14ビットのダイナミック・レンジにスケールされる。

3) Mu-law 圧縮

各音声標準は、図1のブロック186で、参照文献(7)のような公の文献でよく知られているMu-law圧縮を使用して16ビットの精度から8ビットの精度に引き下げられる。8ビットへ圧縮されたデータは、ふたたびMu-law公式に従って、14ビットへ拡大される。

【0034】この結果、図1ブロック188でシミュレートされる電話チャネル音声データ・セットになる。これは、信号強度によって増大、減少する量子化ノイズ・レベルを持ち、およそ一定のS/N比を維持する。特に、話し手の声の大きい場合、これは、電話音声信号の中で聞かされる「ひび割れ」雑音を導入する。

【0035】電話データより高域の種々の帯域幅で集められるであろう音声データ180のこのような処理は、電話機器での使用のため音声認識装置50をブロック190で訓練するために使用される。音響訓練190は、図1のブロック192の音楽モデルと量子化コード・ブック105を生成する。これにより、スフィンクス音声認識エンジンを使用して電話帯域幅での実際の音声認識を行うことを可能とする。

【0036】シミュレートされた電話チャネルデータ使用の認識装置訓練

2個のコード・ブック105と2個の音楽モデル・セット192が作成されるように、2つの訓練セッション、

すなわち電話と高帯域に対するセッションが、実行される。高帯域、ローカルな認識あるいは、電話帯域幅などのユーザの要求に応じてコード・ブック105の各セットおよび各音楽モデル192は、別々にに保管され、実行される。いずれの帯域幅でも、自己相関係数は、ケプストラム係数を引き出すために抜き出される。そのフレームにもっとも近い係数を類別するために、ケプストラム係数がベクトル量子化104によって実行される。このようにして、[]で記述されるように、各音声時系列フレームは、そのフレームを表わす3バイトに渡じられる。

【0037】量子化の値のセットが、ビーム・サーチ・プロセス106に送り出される。ビーム・サーチ106は、ビタービ(Viterbi)ビーム・サーチと呼ばれる文法ガイド型「隠れたマルコフ・モデル」(Hidden Markov Model)サーチ・プロセスである。この文法ガイド型サーチは、サーチ・スペースを減らすために語対文法を使う。

【0038】本発明のもうひとつの重要な点は、その音声認識システムがローカルであろうが遠隔地であろうが、両方の音声処理することができることである。これは、音声のいずれのタイプもチャネル・シミュレータで使われる帯域幅に対応するように、実行時データ条件づけ・速度変換フィルタの遮断ポイントを2個の帯域幅の幅が狭い方に近い帯域幅に置くことによって、達成される。3dbポイントおよび移行帯域特性は、訓練の中で使われる電話符号化器通過帯域の上位移行帯域の特性に近似しなければならない。

【0039】ビーム・サーチ106は、そのベクトル量子化の中で引き出された時系列を語対文法からの語列に突き合わせ、各文脈を定義する。音声認識サーバは、ユーザ適用業務または音声認識クライアント(ブロック110)とコミュニケーションする。本発明の構造は、単数のバックエンドとコミュニケーションする複数のフロントエンド(ワークステーション)または複数のバックエンドとコミュニケーションする複数のフロントエンドを持つことができる。

【0040】本発明のシステムは、オペレーションの異なるレベルのために構成され実行される。非常に高いデータ速度をもつコミュニケーション・ネットワークについては、フロントエンドでのデータ圧縮のために、音声標準は、直接バックエンドを実行しているシステムに伝達されることが可能である。原デジタル音声データストリームが、複数のユーザ用のバックエンドがあるサーバに送り出されることが可能である。電話システムについては、複数のチャネルが1つのバックエンドへつながるか、または、複数のユーザが、フロントエンドおよびバックエンド双方にコミュニケーションする。

【0041】本発明でのシステムは、音声認識サーバとして配備される音声認識機能を中心に主として構成さ

れる。システムは、その時点の文脈として適用業務が選択する語対文法によってガイドされる。音声認識適用業務は、初期値設定プロシージャ、ステータス・コードおよびコマンド[G]のような機能をサポートする適用業務プログラム・インタフェース(API)コールをもつ音声認識システムにインターフェースを持つ。音声認識適用業務は、音声認識サーバに一定のタイプの操作を要求するか、あるいは、ある特定の認識文脈をロードして、必要なこと、音声認識のための文脈を起動するよう要求する。音声認識適用業務が最初に実行されるとき、タスクは通常サーバによって事前ロードされる。適用業務の活動の必要に応じて、タスクはその後順に起動される。

【0042】音声認識サーバ(ブロック108)のAPIコールは、ユーザ適用業務(ブロック110)が音声認識システムのサービスを要求することを可能にする。ユーザ適用業務プログラム(ブロック109)は、音声認識サーバの種々の構成要素と同一コンピュータまたは異なるコンピュータの上で実行することができる。同じコンピュータ上の場合、適用業務プログラム(ブロック110)は、そのオペレーティングシステムでサポートされる共有メモリおよびセマフォを通して音声認識サーバとインターフェースをとることができる。異なるコンピュータ上の場合、通信はRS232Cインターフェースあるいは遠隔プロシージャ呼出し(RPC)を通して行われる。RPCは参照プログラミング文献[10]でよく知られている。

【0043】ユーザ適用業務の典型的例には、エグゼクティブ情報システム、言葉の照会結ぶデータベース・アクセス、ソフトウェア問題報告システムなどがある。【0044】もうひとつの例は、その利点を活用するため音声認識サーバへの呼び出しを行う電話回音音声応答装置(VRU)である。RISC SYSTEM 6000(TM)およびOS/2(TM)をもつPS/2(TM)の上でこれらのサーバは実行された。

【0045】Direct Talk 6000(TM)は、同様の電話VRUシステムである。このVRUシステムでは、1本の電話回線を扱うのではなく、(同時に活動中となる可能性のある24個の会話チャネルをもつ)T1回線処理が必要となる。音声認識サーバ処理は、Direct Talk(TM)のように大量の電話適用業務の処理が必要な場合、複数のクライアントを扱うことができる。ユーザ適用業務は多くの文脈を前もって登録することができる。レストラン案内、ハードディスク・ヘルプ・デスク、あるいは、ソフトウェア・ヘルプ・デスクは全て複数の文脈を階層的に事前に登録することができる。各適用業務では、何人かのユーザが、音声ストリームを入力することができる。各適用業務は、特特有の音声ストリームのために特有の文脈の下で音声認識を実行するよう音声認識サーバに指示する。

【0046】言い換えると、同じAPIを扱う複数のユ

ーザが、1またはおそらくいくつかの版の音声認識サーバを用いるタスクすべてを登録するであろう。システムは、要請された作業がすでにロードされているかを確認し、複数のユーザの音声認識タスクが余分にロードされることを回避する。

【0047】タスク構築(ブロック130)は、いくつかの基本入力ソースを持つ。20,000語の発音をもつ基本辞書である米語辞書(ブロック132)は、その1つである。補足辞書(ブロック138)は、適用業務特有のもので、基本辞書の中で見つけられなかった語の発音を追加するためのものである。補足辞書は、典型的には、特定の適用業務が音声認識のために必要とする固有名詞、頭字語(ACRONYM)その他から構成される。

【0048】基本米語辞書(ブロック132)は、タスク構築プログラム(ブロック134)によって求められる語および音素を供給する。タスク構築プログラムは、また、何がそのタスクの下で音声認識サーバによって認識されることができるかを決めるためにタスクBNF辞書(ブロック136)から該当するタスクBaukus-Naur Form(BNF)文法を引き出す。たとえば、地域レストラン情報を提供する適用業務の最初の文脈は、その話し手が希望するレストランのタイプ、たとえば、フランス、イタリア、中国料理などであるかもしれない。ひとたびそのタイプが決まれば、次の文脈は、その特定のカテゴリの中のレストランとなる。タスク構築プログラムは、そのパターン合わせのために必要なすべての語を見つめるためにBNFを分析し、汎用の米語辞書(ブロック132)から音素表示を引き出す。必然的に、あらゆる特定適用業務は、そのシステムに加えられなければならないそれ自身の副語彙を持ち、それらは、補足辞書に保存される。たとえば、レストラン・ヘルプ・デスクの中で、「イタリアン」、「フレンチ」、「スパニッシュ」などの言葉は、汎用米語辞書で見つけられるが、レストラン名、とくに外国語で、たとえば、「Cherchez Les Femmes」、「Chateau Voulze」や、アメリカのレストランで普通でない名、たとえば、J. J. Muldoon、は、普通の辞書になく、タスク補足辞書(ブロック138)に加えなければならない。これらの補足辞書(ブロック138)は、また、基本汎用米語にあるが発音をローカルなものにするためにローカルな語彙を含めることができる。

【0049】タスク構築プログラム(ブロック134)は、入力BNF文法を分析して、その文法の中の各語のリストと次に続くことができるすべての語のサブリストを生成する。したがって、その文法の中の各語が、後に続く適切な語のリストおよび各語の音素表示のポインタを持つ。音素モデル192は、種々のVQ値を観察するである。このマルコフ・モデルは、VQ値(ブロック104)のための、一群の離散的確率分布であり、「隠れたマルコフ」状態機械が音素の範囲内の特定の状態に

13

あるとすると、VQ値のオカレンスの確率を与える。
「隠れたマルコフ・モデル」は文献[11]に適切に記述されている。

【0050】ビーム・サーチ(ブロック106)は、訓練プロセスの間に生成された文脈認知のトリフォン(tri phones)の大きいテーブルから連結HMM音素モデル192でできている語モデルを使用する。この語モデルが、VQ値の観察された順序を最もよく説明する語順序の最適推定を行うために使われる。ビーム・サーチ(ブロック106)は、そのサーチの中で使われる語をつくるための音素モデル192を選択するために、語文法を使う。

【0051】ユーザ適用業務は、音声認識サーバを制御する。例えば、[12]で記述されるIBMプログラム・プロダクト Direct Talk/2(TM)は、電話に
応答しレストラン案内機能を実行するひとつのユーザ適用業務となり得る。レストラン案内適用業務は、Direct Talk/2(TM)を使用し、この適用業務が16の文脈を持ち、レストランメニュー・デスクの一部である文脈を事前ロードする要求を起こすことを音声認識サーバに知らせる。その適用業務が進行するにつれて、音声認識サーバの文脈切り替えを要請する。ユーザは、電話ヘルプを電話を通して呼び出す。レストラン案内は、音声認識サーバに最初のレベルの文脈での音声認識を実行することを要請する。認識サーバとユーザ適用業務間のAPI上で制御とデータが交換される。Direct Talk/2(TM)システムの複数の事例が同じ音声認識サーバを使用する可能性がある。

【0052】音声認識サーバは、無音間隔(ユーザが調整可能で、ほとんど一般に0.6秒)が来るまで音声データを捕捉する。無音間隔が観察されると、認識は終了し、話し手の話が終わったと仮定される。

【0053】本発明記載の音声認識システムは、複数のハードウェア・プラットフォームおよび複数のソフトウェア機械構成の上に、複数の実施が可能にするよう基本設計がなされる。たとえば、1つの可能な構造は、図5のように、ローカル・エリア・ネットワーク160を通して接続されているワークステーションの物理的実施の上への上記論理的構造50の物理マッピングを提供する。この構造の中の各ワークステーション150、150'、150''、150'''は、複数の独立ユーザ適用業務を実行することができ、各々は、スレーブ・プロセッサとしての音声認識サーバ50のマスターとなる。PBX170は、外部の電話回線に接続して、電話帯域幅データストリームを図1の音声認識サーバ50のアナログ・デジタル変換100に渡す。認識された音声を表わしているテキストが、ワークステーション150、150'、150''、150'''のユーザ適用業務に音声認識サーバから返される。

【0054】訓練プロセス

14

ビーム・サーチ・プロセスでの語モデルとテキストとのパターン合わせに使われる音素HMM192のパラメータを推定するために、訓練プロセスでは、既知の音声およびテキスト原稿という大規模辞書を使用する。

【0055】最初に、その原稿が、訓練セットの語の発音を表わす音素を汎用米語辞書から検索するために使われる。

【0056】次に、音素HMM192のパラメータが、共調音(coarticulation)効果の効果的な推定を行うために、先行および後続音素文脈(トリフォン-tri phonesと呼ばれる)の中で推定される。使われる推定プロセスは、[11]で記述のBaum-Welch 順方向・逆方向繰返しアルゴリズムである。訓練されたトリフォンHMMが訓練セットの中で観察されたVQ値時系列を生成したであろう確率を最大にするために、HMMのパラメータが、繰り返し調節される。

【0057】あらゆる「隠れたマルコフ」音素モデルには多くのパラメータがあり、各「隠れた」状態機械中に7つの状態および12のトランジション・アーク(TRANSITION ARK)が存在する。各トランジション・アークに関連して、3つのコード・ブックの各々の確率分布に、関連する256の離散要素がある。訓練プロセスから生じるトリフォンHMMパラメータは、連続音声の中に存在する共調音効果を表わすのに必要なトリフォン数を減らすために一定の幅の値の範囲に集められる。

【0058】訓練は、ローカルな電話交換を通して集められる低レベル帯域幅音声およびマイクからの高帯域音声の組合せによって実行される。高帯域音声は、本発明に従って、本願書で記述の電話チャネル・シミュレータ185によって処理される。3個のコード・ブックすべては、この段階でコンパイルされる。[1]で記述のように、コード・ブックが、ケプストラム、微分のケプストラム、べきおよび微分のべきを含む。

【0059】3つのコード・ブックの各々は量子化コード・ブック105に保存され、実行時ベクトル量子化プロセスで使われる。ここで、電話ネットワークの効果は、データの事前処理によってシミュレートされ、公衆電話ネットワークが調整されると同じ方法で特性コード・ブックの統計的臨界性が調整されると、このプロセスをとることによって、米国の様々なロケーションからの呼び出しをもつ実際の電話の音声認識の正確度が大幅に増加した。

【0060】図6は、たとえばPBX170経由で電話から得られた音声にตอบสนองする音声認識装置50を訓練するための電話チャネル・シミュレーション・プロセッサ200を記述する流れ図である。図6の流れ図は、図5のデータ処理装置50の上で実行されることができコンピュータ・プログラム方法を表わす。

【0061】プロセス200は、電話帯域幅より帯域幅がより高い音声認識訓練プロセッサ50に音声データ・

15

セットを入力するステップ202で始まる。例となる高帯域音声データ・セットは、参照文献(14)から(19)で記述されている。このステップは、図1のデータ入力ブロック180に対応する。

【0062】図6のステップ204で、音声データ・セットは、電話帯域幅を持つ間引かれた音声データ・セットを得るために間引かれる。これは、図1の機能ブロック182に対応する。間引かれた音声データ・セットは、入力音声データ・セットの高い方の帯域幅より低い帯域幅を持つであろう。間引き(decimation)プロセスは、参照(2)で記述されている。

【0063】次に、図6のステップ206で、帯域通過デジタル濾波器を間引かれた音声データ・セットに適用し、電話機器の伝送特性に特徴づける。これは、図1の機能ブロック182に対応する。これは、帯域された音声データ・セットを得るために行われる。帯域通過デジタル濾波器は、最大平坦設計アルゴリズムを持たなければならない。

【0064】次に、図6の中のステップ208で、その最大ダイナミック・レンジが非圧伸電話音声の最大レンジと一致するように、濾波された音声データ・セットの振幅が、再補正される。これは、図1の機能ブロック184に対応する。これは、振幅再補正音声データ・セットを得るために行われる。このステップの結果、その最大ダイナミック・レンジは非圧伸Muelaw電話音声の最大ダイナミック・レンジと一致し得る。代わりに、その最大ダイナミック・レンジは非圧伸A-law電話音声の最大ダイナミック・レンジと一致することもできる。

【0065】次に、図6のステップ210で、上記補正音声データ・セットを、電話中の音声信号の圧伸非圧伸の順序を表わしている量子化ノイズをもって修正する。これは、図1の機能ブロック186に対応する。これは、修正された音声データ・セットを得るために行われる。修正ステップは、Muelawノイズとしての量子化ノイズを持つことができる。代わりに、修正ステップは、A-lawノイズとしての量子化ノイズを持つことができる。

【0066】次に、図6のステップ212では、統計的パターン・マッチング・データ装置を訓練するために、音声認識プロセッサ50へ修正された音声データ・セットを入力する。これは、図1の出力データ・ブロック188に対応する。シミュレートされた電話チャネル音声185が、電話音声特有性を持つ電話コード・ブック105の特性を持つ音素モデル192を生成するために、音響的訓練プロセッサ190によって使われる。

【0067】次に、図6のステップ214で、たとえば、図5のPBX170からの信号のような、電話からの音声信号に対し、音声認識プロセッサ50を使って、音声認識が実行される。

16

【0068】電話チャネル・シミュレータ(ブロック185)を使用する高帯域音声の変換は、連続の音声認識装置に限られてなく、たとえば、IBM Tangora Dictation SystemおよびDragon Systems、ニュートン・マサチューセッツ、Dragon 30K DictateおよびKurzwel Applied Intelligence, Voice Report, Waltham、マサチューセッツおよび(20)で記述されるその他のシステム等のような様々な音声認識プロセッサに適用されるということに留意する必要がある。

【0069】上記本発明の説明において引用した参照文献は、以下の通りである。

【0070】(1) "Large Vocabulary Speaker and Dependent Continuous Speech Recognition: The Sphinx System"; Kai-Fu Lee; Carnegie Mellon University, Department of Electrical and Computer Engineering; April 1988; CMU-CS-88-148

(2) "A General Program to Perform Sampling Rate Conversion of Data by Rational Ratios"; from "Programs for Digital Signal Processing", Ed.: Digital Signal Processing Committee of the IEEE Acoustics, Speech, and Signal Processing Society; IEEE Press, 1979; Section 8.2, pp.8.2-1 to 8.2-7 by R.E. Crochiere

(3) "Theory and Application of Digital Signal Processing" L.R. Rabiner, B. Gold; Prentice Hall, 1975, pp 91

(4) "Digital Processing of Speech Signals"; Prentice Hall Signal Processing Series; 1978, pp 401-402, 411-413

(5) "An Algorithm for Vector Quantizer Design"; Y. Linde, A. Buzo, R. Gray, IEEE Transactions on Communications, Vol. com-28, no. 1, January 1980

(6) "IBM Continuous Speech Recognition System Programmers Guide"; B. Booth; 1992; currently unpublished, available on request.

(7) "Digital Telephony and Network Integration"; B. Keiser, E. Strange; Van Nostrand Reinhold Company Inc. 1985.; pp. 26-31

(8) "Design Subroutine (MAXFLAT) for Symmetric FIR Low Pass Digital Filters with Maximally-Flat Pass and Stop Bands" from "Programs for Digital Signal Processing", Ed.: Digital Signal Processing Committee of the IEEE Acoustics, Speech, and Signal Processing Society; IEEE Press, 1979; Section 5.3, pp 5.3-1 to 5.3-6 by J. Kaiser

(9) "Acoustical and Environmental Robustness in Automatic Speech Recognition" A. Acero; Carnegie Mellon University, Department of Electrical and Computer Engineering; April 1990; CMU-CS-88-148

(10) "AIX Distributed Environments: NFS, NCS, RPC,

DS Migration, LAN Maintenance and Everything"; IBM International Technical Support Centers, Publication GG24-3489, May 8, 1990

[11] "A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition"; L. Rabiner; Readings in Speech Recognition; Ed.: A. Waibel, K. Lee; Morgan Kaufmann; 1990; pp 267-296

[12] "IBM CallPath DirectTalk/2 General Information and Planning Manual"; International Business Machines publication no. GB35-4403-0; 1991

[13] "A Maximum Likelihood Approach to Continuous Speech Recognition"; L. R. Bahl, F. Jelinek, R. Mercer; Readings in Speech Recognition; Ed.: A. Waibel, K. Lee; Morgan Kaufmann; 1990; pp 308-319

[14] "Speech Corpora Produced on CD-ROM Media by The National Institute of Standards and Technology (NIST)", April, 1991

[15] "DARPA Resource Management Continuous Speech Database (RMI) Speaker Dependent Training Data", September 1989 NIST Speech Discs 2-1.1, 2-2.1 (2 Discs) NTIS Order No. PB89-226666

[16] "DARPA Resource Management Continuous Speech Database (RMI) Speaker-Independent Training Data", November 1989 NIST Speech Disc 2-3.1 (1 Disc) NTIS Order No. PB90-500539

[17] "DARPA Extended Resource Management Continuous Speech Speaker-Dependent Corpus (RM2)", September 1990 NIST Speech Discs 3-1.2, 3-2.2 NTIS Order No. PB90-501776

[18] "DARPA Acoustic-Phonetic Continuous Speech Corpus (TIMIT)", October 1990 NIST Speech Disc 1-1.1 NTIS Order No. PB91-0505065

[19] "Studio Quality Speaker-Independent Connected-Digit Corpus (TIDIGITS)", NIST Speech Discs 4-1.1, 4-2.1, 4-3.1 NTIS Order No. PB91-505592

[20] "The Spoken Word", Kai-Fu Lee, et al., Byte Magazine, July 1990, Vol- 15, No. 7; pp. 225-232 [0071]

【発明の効果】電話回線から入力される不特定の話し手の音声による音声認識システムを構築することによって、たとえば、電話による顧客問い合わせ自動応答システムやレストラン電話案内など、従来技法では実現できなかった新たなコンピュータ適用業務を開発することができる。

【図面の簡単な説明】

【図1】 電話チャネル・シミュレータ発明を含む、連続音声認識システムの論理的構造を図示する。

10 【図2】 電話の符号器器渡波器インパルス応答を特徴づけるグラフである。

【図3】 振幅特性額文規格化ラジアン周波数を図示するグラフである。

【図4】 対数振幅特性額文規格化ラジアン周波数を図示するグラフである。

【図5】 電話顧客業務通話センタにおける音声認識サーバのネットワークのブロック図である。

【図6】 電話から得られる音声に 대응するために音声認識装置を訓練するためのプロセスのステップ流れ図である。

20 【符号の説明】

100 アナログ・デジタル変換

104 ベクトル量子化

105 ベクトル量子化コードブック

192 音楽モデル

135 語対文法

132 米語辞書

138 補助辞書

186 Mue-lawノイズ

186 A-lawノイズ

188 電話チャネル・シミュレータ

182 符号器デジタル渡波・速度変換

184 振幅補正 (スケージング)

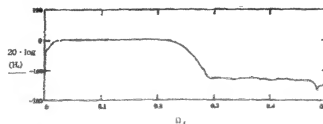
134 タスク構築プログラム

106 ビーム・サーチ

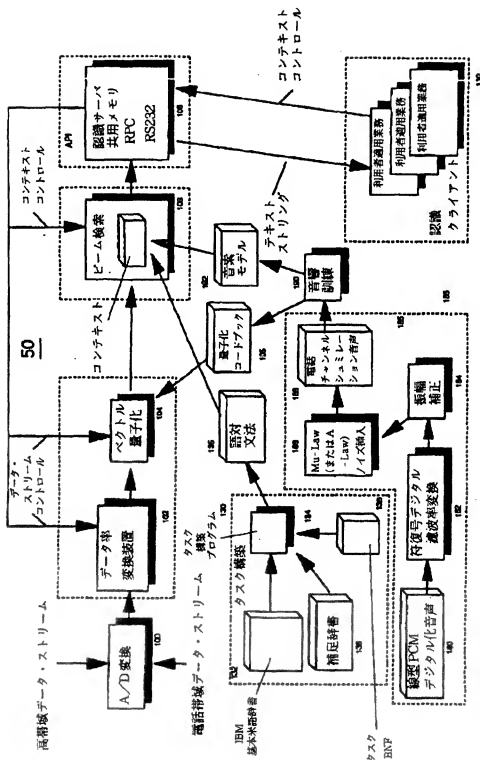
108 API (適用業務プログラム・インターフェー

ス)

【図4】



【図1】



ICSPS システム論理構造

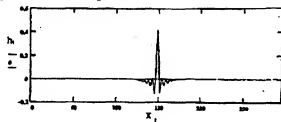
【図2】

N := READ (FilterSizeFile) N = 255

t := 0..N-1

x_t := t

h_t := READ (FilterWeightFile)



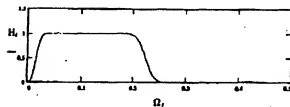
【図3】

F := 256

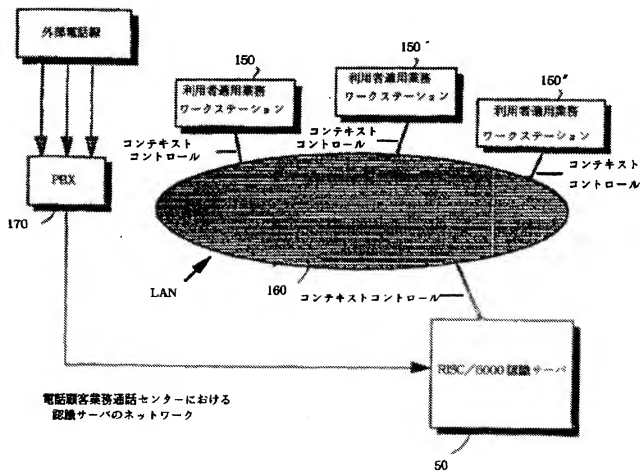
f := 0..F

$\Omega_f := \frac{f}{2F}$

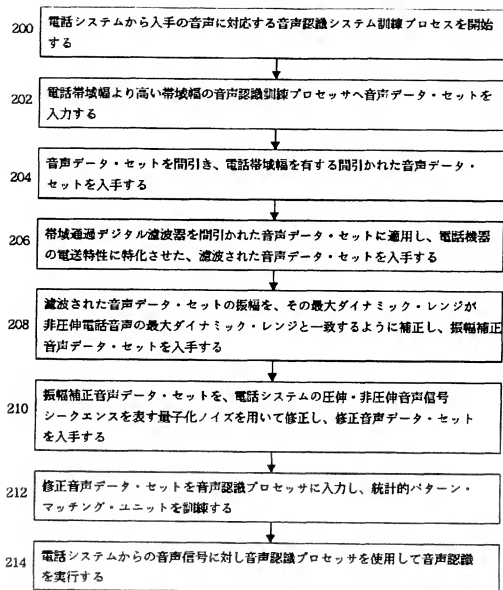
$$H_f := \left| \sum_t h_t \cdot e^{-i \cdot 2 \cdot \pi \cdot t \cdot \frac{f}{2 \cdot F}} \right|$$



【図5】



【図6】



フロントページの続き

(72)発明者 ノーマン エフ ブリックマン
 アメリカ合衆国メリーランド州ボトマック
 ミルバン ドライブ 11709番地